

METHODS AND APPARATUS FOR OUTLIER DETECTION FOR HIGH DIMENSIONAL DATA SETS

Field of the Invention

The present invention is related to outlier detection in high dimensional data and, more particularly, to methods and apparatus for performing such detection in accordance with various high dimensional data domain applications where it is important to be able to find and detect outliers which deviate considerably from the rest of the data.

Background of the Invention

The outlier detection problem is an important one for very high dimensional data sets. Much of the recent work has focused on finding outliers for high dimensional data sets which are based on relatively low dimensionalities, for example, up to 10 or 20. However, the typical applications in which points are outliers may involve higher dimensionality such as, for example, 100 or 200. For such applications, more effective techniques are required for outlier detection.

Many data mining algorithms described in the literature find outliers as an aside-product of clustering algorithms. Such techniques typically find outliers based on their nuisance value rather than using techniques which are focused towards detecting deviations, see, e.g., A. Arning et al., "A Linear Method for Deviation Detection in Large Databases," Proceedings of the KDD Conference, 1995. Outliers are however quite useful based on their value for finding behavior which deviates significantly from the norm. In this invention, we carefully distinguish between the two, and develop algorithms which generate only outliers which are based on their deviation value.

Although the outlier detection definition described in S. Ramaswamy et al., "Efficient Algorithms for Mining Outliers from Large Data Sets," Proceedings of the ACM SIGMOD Conference, 2000 has some advantages over that provided in E. Knorr et al., "Algorithms for Mining Distance-based Outliers in Large Data Sets," Proceedings of the VLDB Conference, September 1998, both of them suffer from the same inherent

5

10

0
9
8
7
6
5
4
3
2
15
20

disadvantages of treating the entire data in a uniform way. However, different localities of the data may contain clusters of varying density. Consequently, a new technique which finds outliers based on their local density was proposed in M.M. Breunig et al., “LOF: Identifying Density-Based Local Outliers,” Proceedings of the ACM SIGMOD Conference, 2000, which finds the outliers based on their local neighborhoods; particularly with respect to the densities of these neighborhoods. This technique has some advantages in accounting for local levels of skews and abnormalities in data collections. In order to compute the outlier factor of a point, the method in the M.M. Breunig et al. reference computes the local reachability density of a point o by using the average smoothed distances to a certain number of points in the locality of o .

25

Thus, the above-mentioned techniques proposed in the above-cited E. Norr et al. reference, the S. Ramaswamy et al. reference and the M.M. Breunig et al. reference try to define outliers based on the distances in full dimensional space in one way or another. Recent results have also shown that when the distances between pairs of points are measured in the full dimensional space, all pairs of points are almost equidistant, see, e.g., K. Beyer et al., “When is Nearest Neighbors Meaningful?” Proceedings of the ICDT, 1999. In such cases, it becomes difficult to use these measures effectively, since it is no longer clear whether or not these are meaningful. In the context of the algorithms proposed in the above-cited E. Knorr et al. reference, a very small variation in d can result in either all points being considered outliers or no point being considered an outlier. The definition in the S. Ramaswamy et al. reference is slightly more stable since it does not rely on the use of such a parameter which is difficult to pick a priori. However, for high dimensional problems, the meaningfulness of the k -nearest neighbor in high dimensional space is in itself in doubt; therefore, the quality of outliers picked by such a method may be difficult to estimate. The same problem is relevant for the method discussed in the M.M. Breunig et al. reference in a more subtle way; since the local densities are defined using full dimensional distance measures.

For problems such as clustering, it has been shown (e.g., in C.C. Aggarwal et al., "Fast Algorithms for Projected Clustering," Proceedings of the ACM SIGMOD Conference, 1999 and C.C. Aggarwal et al., "Finding Generalized Projected Clusters in High Dimensional Spaces," Proceedings of the ACM SIGMOD Conference, 2000) that by examining the behavior of the data in subspaces, it is possible to design more meaningful clusters which are specific to the particular subspace in question. This is because different localities of the data are dense with respect to different subsets of attributes. By defining clusters which are specific to particular projections of the data, it is possible to design more effective techniques for finding clusters. The same insight is true for outliers, because in typical applications such as credit card fraud, only the subset of the attributes which are actually affected by the abnormality of the activity are likely to be applicable in detecting the behavior.

In order to more fully explain this point, let us consider the example illustrated in FIGs. 1A-1D. In the example, we have shown several 2-dimensional cross-sections of a very high dimensional data set. It is quite likely that for high dimensional data, many of the cross-sections may be structured; whereas others may be more noisy. For example, the points A and B show abnormal behavior in views 1 (FIG. 1A) and 4 (FIG. 1D) of the data. In other views, i.e., views 2 (FIG. 1B) and 3 (FIG. 1C), the points show average behavior. In the context of a credit card fraud application, both the points A and B may correspond to different kinds of fraudulent behavior, yet may show average behavior when distances are measured in all the dimensions. Thus, by using full dimensional distance measures, it would be more difficult to detect points which are outliers, because of the averaging behavior of the noisy and irrelevant dimensions. Furthermore, it is impossible to prune off specific features *a priori*, since different points (such as A and B) may show different kinds of abnormal patterns, each of which use different features or views.

Thus, the problem of outlier detection becomes increasingly difficult for very high dimensional data sets, just as any of the other problems in the literature such as

clustering, indexing, classification, or similarity search. Previous work on outlier detection has not focused on the high dimensionality aspect of outlier detection, and has used methods which are more applicable for low dimensional problems by using relatively straightforward proximity measures, e.g., the above-mentioned E. Knorr et al. and S. Ramaswamy et al. references. This is very important for practical data mining applications which are mostly likely to arise in the context of very large numbers of features. The present invention focuses for the first time on the effects of high dimensionality on the problem of outlier detection. Recent work has discussed some of the concepts of defining the intentional knowledge which characterizes distance-based outliers in terms of subsets of attributes. Unfortunately, this technique was not intended for high dimensional data, and the complexity increases exponentially with dimensionality. As the results in E. Knorr et al., "Finding Intentional Knowledge of Distance-based Outliers," Proceedings of the VLDB Conference, September, 1999 show, even for relatively small dimensionalities of 8 to 10, the technique is highly computationally intensive. For even slightly higher dimensionalities, the technique is likely to be infeasible from a computational standpoint.

Summary of the Invention

The present invention provides methods and apparatus for outlier detection which find outliers by observing the density distributions of projections from the data. Intuitively, this new definition considers a point to be an outlier, if in some lower dimensional projection, it is present in a local region of abnormally low density. Specifically, the invention defines outliers for data by looking at those projections of the data which have abnormally low density.

Accordingly, in an illustrative aspect of the invention, a method of detecting one or more outliers in a data set comprises the following steps. First, one or more sets of dimensions and corresponding ranges (e.g., patterns) in the data set which are sparse in density (e.g., have an abnormally low presence which cannot be justified by randomness)

are determined. Then, one or more data points (e.g., records) in the data set which contain these sets of dimensions and corresponding ranges are determined, the one or more data points being identified as the one or more outliers in the data set.

In further illustrative aspects of the invention, the range may be defined as a set of contiguous values on a given dimension. The sets of dimensions and corresponding ranges in which the data is sparse in density may be quantified by a sparsity coefficient measure. The sparsity coefficient measure $S(D)$ may be defined as $\frac{n(D) - N * f^k}{\sqrt{N * f^k * (1-f^k)}}$, where k represents the number of dimensions in the data set, f represents the fraction of data points in each range, N is the total number of data points in the data set, and $n(D)$ is the number of data points in a set of dimensions D . A given sparsity coefficient measure is preferably inversely proportional to the number of data points in a given set of dimensions and corresponding ranges. A set of dimensions may be determined using an algorithm which uses the processes of solution recombination, selection and mutation over a population of multiple solutions. The process of solution recombination may comprise combining characteristics of two solutions in order to create two new solutions. The process of mutation may comprise changing a particular characteristic of a solution in order to result in a new solution. The process of selection may comprise biasing the population in order to favor solutions which are more optimum.

These and other objects, features and advantages of the present invention will become apparent from the following detailed description of illustrative embodiments thereof, which is to be read in connection with the accompanying drawings.

Brief Description of the Drawings

FIGs. 1A-1D are diagrams of various patterns of data sets illustrating outlier detection issues;

FIG. 2 is a block diagram illustrating a hardware implementation suitable for employing outlier detection methodologies according to an embodiment of the present invention;

FIG. 3 is a flow diagram illustrating an overall process for outlier detection according to an embodiment of the present invention;

FIG. 4 is a flow diagram illustrating a procedure for encoding potential solutions as strings according to an embodiment of the present invention;

FIG. 5 is a flow diagram illustrating a procedure for selection used by a genetic outlier detection algorithm according to an embodiment of the present invention;

FIG. 6 is a flow diagram illustrating a procedure for crossover and solution recombination used by a genetic outlier detection algorithm according to an embodiment of the present invention;

FIG. 7 is a flow diagram illustrating a procedure for mutation used by a genetic outlier detection algorithm according to an embodiment of the present invention; and

FIG. 8 is a diagram illustrating a broad outline of how the multi-population hill climbing, recombination and search space exploration works.

Detailed Description of Preferred Embodiments

As mentioned above, the present invention provides a new technique for outlier detection which finds outliers by observing the density distributions of projections from the data. Intuitively, this new definition considers a point to be an outlier, if in some lower dimensional projection, it is present in a local region of abnormally low density.

Specifically, the invention defines outliers for data by looking at those projections of the data which have abnormally low density. Thus, a first step is to identify and mine those patterns which have abnormally low presence which cannot be justified by randomness. Once such patterns have been identified, then the outliers are defined as those records which have such abnormal patterns present in them.

5

In order to find such abnormal lower dimensional projections, we need to define and characterize what we mean by an abnormal lower dimensional projection. An abnormal lower dimensional projection is one in which the density of the data is exceptionally lower than average. In order to find such projections, we first perform a grid discretization of the data. Each attribute of the data is divided into p ranges. These ranges are created on an equi-depth basis; thus each range contains a fraction $f = 1/p$ of the records. These ranges form the units of locality which we will use in order to define low dimensional projections which have unreasonably sparse regions.

10

Let us consider a k -dimensional region. The expected fraction of the records in that region, if the attributes were statistically independent, would be equal to f^k . Of course, the data is far from statistically independent; it is precisely the deviations which are abnormally below the norm which are useful for the purpose of outlier detection.

15

Let us assume that there are a total of N points in the database. Then, the expected fraction and standard deviation of the points in a k -dimensional cube is given by $N * f^k$ and a standard deviation of $\sqrt{N * f^k * (1 - f^k)}$. Let $n(D)$ be the number of points in a k -dimensional cube D . Then, we calculate the sparsity coefficient $S(D)$ of the cube D as follows:

$$S(D) = \frac{n(D) - N * f^k}{\sqrt{N * f^k * (1 - f^k)}}$$

20

Only sparsity coefficients which are negative indicate cubes which have lower presence than expected. Thus, it is desirable to find those projections which have low (or highly negative) sparsity coefficients.

25

Now we will discuss the algorithms which are useful for outlier detection in high dimensional problems. A natural class of methods for outlier detection are the naive brute-force techniques in which all subsets of dimensions are examined for possible patterns which are sparse. These patterns are then used in order to determine the points

which are possibly outliers. We propose a naive brute-force algorithm which is very slow at finding the best patterns because of its exhaustive search of the entire space, and a much faster algorithm which is able to quickly prune away large parts of the search space.

The problem of finding subsets of dimensions which are sparsely populated is a difficult one, since one needs to look at an exponential number of combinations of attributes in order to find outliers. Furthermore, it may often be the case that even though particular regions may be well populated on certain sets of dimensions, they may be very sparsely populated when such dimensions are combined together. For example, there may be a large number of people below the age of 20, and a large number of people with diabetes; but very few with both. Consequently, it becomes difficult to prune the search space using structured search methods. Consequently, we borrow techniques from the class of genetic algorithms in order to design effective techniques for high dimensional outlier detection. For this purpose, we carefully design the various genetic algorithm components by effectively designing the data encoding, crossover and mutation techniques which are structurally suited to the outlier detection problem.

The idea of using the principles of organic evolution for solving combinatorial optimization problems was introduced by John Holland about thirty years ago. This idea was subsequently formalized by him in 1975 as Genetic Algorithms (D.E. Goldberg, "Genetic algorithms in search, optimization and machine learning," Addison Wesley, Reading, MA, 1989). In his seminal work, Holland laid down the theoretical foundations of the area, and paved the way for all subsequent Genetic Algorithm research. In the past decade, the field of Genetic Algorithms has seen rapid progress both in terms of theoretical as well as applied work.

Genetic Algorithms are methods which imitate the process of organic evolution in order to solve parameter optimization problems. The principles of organic evolution were laid down by Charles Darwin several decades ago. The fundamental idea underlying the Darwinian view of evolution is that, in nature, resources are scarce and this automatically leads to a competition among the various species. As a result, all the

species undergo a selection mechanism, in which only the fittest survive. Consequently, the fitter individuals tend to mate with each other more often, resulting in still better individuals. At the same time, once in a while, nature also throws in a variant by the process of mutation, so as to ensure a sufficient amount of diversity among the species, and hence also a greater scope for improvement. The basic idea behind Genetic Algorithms is also similar; every solution to an optimization problem can be “disguised” as an individual in an evolutionary system. The measure of fitness of this “individual” is simply equal to the objective function value of the corresponding solution, and the other species which this individual has to compete with are simply a group of other solutions to the problems. This is one of the reasons why Genetic Algorithms are more effective as heuristic search methods than either hill-climbing, random search or simulated annealing techniques; they use the essence of the techniques of all these methods in conjunction with recombination of multiple solutions in a population. Genetic Algorithms work not with one solution, but with a whole set of them at a time. Appropriate operations are defined in order to imitate the recombination and mutation processes as well, and the simulation is complete. A broad outline of how the multi-population hill climbing, recombination and search space exploration actually works is illustrated in FIG. 8.

Genetic Algorithms have become increasingly important in the past few years as compared to traditional optimization methods. This is primarily because there are large classes of optimization problems for which no efficient algorithms have been developed. Such problems may have an exponential search space, and the distance function may be very noisy and multi-modal, which results in a parameter optimization problem that is treacherously difficult to solve. Many of these problems arise in actual practical situations and require only specifications of approximately optimal solutions rather than provably optimal ones. In such situations, Genetic Algorithms certainly provide an empirically efficient method and perform much better than other traditional approaches such as hill climbing methods.

One of the interesting aspects of Genetic Algorithms is that for every problem the basic Genetic Algorithm used is the same. The only aspect in which the Genetic Algorithm for two problems differ is in the method by which feasible solutions to the combinatorial problem are disguised (or coded) as individuals in the population. Thus, in some sense, the complexity of any problem is captured in the concise problem of representing every solution as an individual in the population in such a way that encourages biological evolution.

Whenever the Genetic Algorithm is used to solve a particular problem, each feasible solution to that problem is defined as an individual. This feasible solution is in the form of a string and is the genetic representation of the individual. Such a string is referred to as a chromosome. In this invention, we will consistently refer to it as a string. Thus, in order to give a genetic representation to our individual, we must have a procedure which converts feasible solutions of the problem into strings which the Genetic Algorithm (hereinafter GA) can recognize and work with. This process of conversion is called coding. For example, in our invention, the string representation of an individual contains the set of dimensions in the record which are included in the projection. The measure of fitness of an individual is evaluated by the fitness function, which has as its argument the string representation of the individual and returns a non-negative real number indicating the fitness value. The fitness value of an individual is analogous to the objective function value; the better the objective function value, the larger the fitness value. Thus, GAs are naturally defined as maximization problems over non-negative objective function values. However, minimization problems can be easily converted into maximization problems on the Genetic Algorithm by simple fitness function transformations.

As the process of evolution progresses, all the individuals in the population tend to genetically become more and more similar to each other. This phenomenon is referred to as convergence. A different method is to terminate the algorithm after a pre-specified number of generations.

We now discuss a genetic algorithm for outlier detection which is able to find outliers by searching for subsets of dimensions in which the data is populated very sparsely. Genetic Algorithmic Techniques, as described in D.E. Goldberg, "Genetic algorithms in search, optimization and machine learning," Addison Wesley, Reading, MA, 1989, are heuristic search methods which rely on successive solution recombinations, random explorations, and selections in order to gradually evolve the most optimum characteristics of a given solution. Problems which are inherently either computationally intensive because of an exponential search space or non-linear/unstructured with respect to the optimization function are good candidates.

However, the exact quality of performance of a genetic algorithm is often dependent on how well it is tailored to a given problems. Typically, genetic algorithms which are customized for given problems in terms of the methods for solution recombination and random explorations perform significantly better than using black-box Genetic Algorithm software on straightforward string encodings, see, e.g., C.C. Aggarwal et al., "Optimized Crossover for the Independent set problem," Operations Research, March 1997. In accordance with the present invention, we provide a genetic algorithm which works effectively for the outlier detection problem.

We now discuss the application of the search technique to the outlier detection problem. Let us assume that the grid range for the i^{th} dimension is denoted by $t(i)$. Then, the value of $t(i)$ can take on any of the values 1 through p , or it can take on the value *, which denotes a "don't care." Thus, there are a total of $p+1$ values that the dimension $t(i)$ can take on. Thus, consider a 4-dimensional problem with $p=10$. Then, one possible example of a solution to the problem is given by *3*9. In this case, the ranges for the second and fourth dimension are identified, whereas the first and third are left as "don't cares." The fitness for the corresponding solution may be computed using the sparsity coefficient provided above.

The genetic algorithm has three main processes; those of selection, crossover and mutation which are performed repeatedly in order to find the interesting projections in

which the outliers exist. We will now discuss the details of these operations in the remainder of detailed description of preferred embodiments below.

The following portion of the detailed description will illustrate the invention using an exemplary data processing system architecture. It should be understood, however, that the invention is not limited to use with any particular system architecture or application. The invention is instead more generally applicable to any data processing system or application in which it is desirable to perform more meaningful outlier detection by observing the density distributions of projections from the data.

FIG. 2 is a block diagram illustrating a hardware implementation suitable for employing outlier detection methodologies according to an embodiment of the present invention. As illustrated, an exemplary system comprises client devices 10 coupled, via a large network 20, to a server 30. The server 30 may comprise a central processing unit (CPU) 32, coupled to a main memory 34 and a disk 36. The main memory 34 may also comprise a cache 38 in order to speed up calculations. It is assumed that multiple clients 10 can interact with the server 30 over the large network 20. It is to be appreciated that the network 20 may be a public information network such as, for example, the Internet or world wide web, however, the clients and server may alternatively be connected via a private network, a local area network, or some other suitable network.

It is assumed that the server 30 contains a large repository of data which is used for the purpose of data mining. The requests for finding the outliers along with the corresponding data sets are specified at the client end 10. These requests are then responded to using the methodologies of the present invention as implemented on the server end 30. The computation is performed by the CPU 32. The data on which the analysis is carried out may already be available at the server on its disk 36, or it may be specified by the client. In either case, the computation is performed at the server end, and the results are returned to and presented to (e.g., displayed) the client.

In one preferred embodiment, software components including instructions or code for performing the methodologies of the invention, as described herein, may be stored in

one or more memory devices described above with respect to the server and, when ready to be utilized, loaded in part or in whole and executed by the CPU.

FIG. 3 is a flow diagram illustrating an overall process for outlier detection according to an embodiment of the present invention. Specifically, the flow diagram of FIG. 3 describes the steps of first encoding the transactions as strings, and then running the iterative genetic algorithm process on the strings in order to find the appropriate outlier projections. These projections along with the corresponding outlier points are returned by the algorithm. The inputs to the process include the projection dimensionality, the set of database points and the number of patterns m . The process starts at block 300. In step 310, an encoding is determined for the database by creating the intervals for each dimension. For example, let us consider a 2-dimensional database in which there are two attributes, age and salary. Then, the encoding will be created by a string of length 2. Let us say that each of the two attributes is divided into $p = 3$ ranges. For example:

Age: Range 1 → 0-30
Range 2 → 31-60
Range 3 → 61 and above
Salary: Range 1 → 0-50,000
Range 2 → 50,001-100,000
Range 3 → 100,001 and above

Then, a 28-year-old person with a salary of 60,000 would be encoded as 12 (i.e., Age Range 1 followed by Salary Range 2), whereas a 62-year-old person with a salary of 20,000 would be encoded as 31. Also, in step 310, we first divide each of the attributes of the database into intervals. Each of these intervals is chosen in such a way that an equal number of records satisfy them. Thus, if p intervals are chosen, then exactly a fraction $1/p$ of the records in the database lie in each of these intervals. In step 320, each record in the database is expressed in terms of these intervals.

In the next set of steps (330-370), an attempt is made to discover those subpatterns of these string representations of the database which are very sparsely represented in the database by using a genetic algorithm-based technique. At this stage, we mention how the genetic algorithm representation technique represents a "don't care"-*. Thus, for a d -dimensional database, if we want to find k -dimensional projections which are sparse, then exactly $(d-k)$ entries of the string would be *, whereas other entries would be a number between 1 through p . For example, in the 2-dimensional example enumerated above, the string $2*$ refers to any individual whose age-range is in 31-60, but the salary could be in any range. This is an example of a 1-dimensional projection of a 2-dimensional problem.

In step 330, a random population P of string solutions to the problem is initialized. Each member of this population is defined by first randomly picking $(d-k)$ positions and setting them to *. The remaining k positions are set to any number between 1 and p . Once the population has been initialized, we run the three operators of selection, mutation and crossover on the strings. The process of selection is performed in step 340, in which those strings which represent more sparse patterns are given greater representation. A detailed description of this process is provided in FIG. 5. In step 350, we perform crossover and recombination of the solution strings in the population. A detailed description of this process is provided in FIG. 6. In step 360, we perform the mutation operation. A detailed description of this process is provided in FIG. 7.

In step 365, the fitness of each string in the population is calculated. Note that the fitness of a string is determined by the number of records in the database that the string covers. The smaller the number of records it covers, the fitter the string, since it is more likely to be a sparsely populated projection. An output list is then the m best strings found so far. In step 370, it is tested whether the process should be terminated at that point. In order to perform the test, several criteria are possible. One of the criteria is whether the looping steps embodied by blocks 340, 350, and 360 have been executed more than a certain number of times. If not, the process returns to step 340. If yes, in

step 380, the set of points in the database which are covered by any of the strings on the output list is found. A point in the database is covered by a string if that point lies in the database projection which is determined by that string. For the 2-dimensional example illustrated above, the database point 21 is covered by the string 2*.

5 In step 390, the sparsely populated projections for the string encodings together with the corresponding database points are reported. The process stops at block 395.

In FIG. 4, we illustrate a process for determining the encoding for each of the records in the database according to an embodiment of the invention. This process is required in order to implement step 320 of FIG. 3. The process starts at block 400. In
10 step 410, the range for each dimension is divided into p equi-depth intervals. Note that in equi-depth intervals, each member of the population is covered by an equal number of records. In step 420, the range corresponding to the i^{th} attribute is found. Thus, for a given point, the i^{th} attribute is made equal to the number of the range corresponding to it. In step 430, the string encoded representation of the object is then returned. The process stops at block 440.
0115

FIG. 5 is a flow diagram illustrating a procedure for selection used by a genetic outlier detection algorithm according to an embodiment of the present invention. This selection process corresponds to step 340 in FIG. 3. The motivation behind performing the selection process is to bias the population in order to make it contain a disproportionately high number of strings which are fit. The process starts at block 500. In step 510, the objective function value of each string solution is computed. The objective function value of a string solution is determined by the number of members of the population which cover that string. This is accomplished by using the sparsity coefficient $S(D)$ described above. The lower this number, the more fit the string. In step
20 520, a ranking of the strings in the population is created, so that the fitter members (i.e., having higher sparsity coefficients) are ranked first, and the least fit members (i.e., having lower sparsity coefficients) are ranked last. The only exception is that we want the string
25

to cover at least one record in the database. Strings which do not cover any member of the population are ranked last.

In step 530, a predetermined fraction of the strings which are lowest in the ranking are removed, and replaced with strings which are higher up in the ranking. This replacement could be done randomly. For example, a random sample could be drawn out of those strings which are higher in the ranking and could be used in order to replace the strings which are lower in the ranking. Thus, step 530 results in a population of strings which correspond to more fitter solutions. In step 540, the biased population P is returned. The process stops at block 550.

FIG. 6 is a flow diagram illustrating a procedure for crossover and solution recombination used by a genetic outlier detection algorithm according to an embodiment of the present invention. This crossover and solution recombination process corresponds to step 350 in FIG. 3. The process starts at block 600. The input to the process is the population before crossover. In step 610, the strings are randomly matched pairwise in the population P. In steps 620 through 640, a looping structure is set up which performs the solution recombination over these pairwise assignments. In order to actually perform the recombination, the strings are iteratively looped over and, for each pair of strings, those positions in which both strings are * are identified. For each such position, the values in the two strings are independently exchanged with a probability of 0.5. This process is performed in step 630, whereas the loop structure is implemented by the steps 620 and 640. In step 645, the updated population is reported, which is created by these recombination operations.

FIG. 7 is a flow diagram illustrating a procedure for mutation used by a genetic outlier detection algorithm according to an embodiment of the present invention. This mutation process corresponds to step 360 in FIG. 3. The process starts at block 700. In order to perform the mutations, a fraction or probability f of the strings in the population is chosen in step 710. The value of f is a user-defined parameter. Next, in step 720, for each of the selected strings, an attribute value which is not * is selected, and mutated to a

random number between 1 and p . In step 730, the set of mutated strings in the population is returned. The process ends at block 740.

Accordingly, as described above in accordance with the present invention, methods and apparatus are provided for outlier detection in databases by determining sparse low dimensional projections. These sparse projections are used for the purpose of determining which points are outliers. The methodologies of the invention are very relevant in providing a novel definition of exceptions or outliers for the high dimensional domain of data.

Although illustrative embodiments of the present invention have been described herein with reference to the accompanying drawings, it is to be understood that the invention is not limited to those precise embodiments, and that various other changes and modifications may be affected therein by one skilled in the art without departing from the scope or spirit of the invention.

10

00
01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
259
260
261
262
263
264
265
266
267
268
269
269
270
271
272
273
274
275
276
277
278
279
279
280
281
282
283
284
285
286
287
288
289
289
290
291
292
293
294
295
296
297
298
299
299
300
301
302
303
304
305
306
307
308
309
309
310
311
312
313
314
315
316
317
318
319
319
320
321
322
323
324
325
326
327
328
329
329
330
331
332
333
334
335
336
337
338
339
339
340
341
342
343
344
345
346
347
348
349
349
350
351
352
353
354
355
356
357
358
359
359
360
361
362
363
364
365
366
367
368
369
369
370
371
372
373
374
375
376
377
378
379
379
380
381
382
383
384
385
386
387
388
389
389
390
391
392
393
394
395
396
397
398
399
399
400
401
402
403
404
405
406
407
408
409
409
410
411
412
413
414
415
416
417
418
419
419
420
421
422
423
424
425
426
427
428
429
429
430
431
432
433
434
435
436
437
438
439
439
440
441
442
443
444
445
446
447
448
449
449
450
451
452
453
454
455
456
457
458
459
459
460
461
462
463
464
465
466
467
468
469
469
470
471
472
473
474
475
476
477
478
479
479
480
481
482
483
484
485
486
487
488
489
489
490
491
492
493
494
495
496
497
498
499
499
500
501
502
503
504
505
506
507
508
509
509
510
511
512
513
514
515
516
517
518
519
519
520
521
522
523
524
525
526
527
528
529
529
530
531
532
533
534
535
536
537
538
539
539
540
541
542
543
544
545
546
547
548
549
549
550
551
552
553
554
555
556
557
558
559
559
560
561
562
563
564
565
566
567
568
569
569
570
571
572
573
574
575
576
577
578
579
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
598
599
599
600
601
602
603
604
605
606
607
608
609
609
610
611
612
613
614
615
616
617
618
619
619
620
621
622
623
624
625
626
627
628
629
629
630
631
632
633
634
635
636
637
638
639
639
640
641
642
643
644
645
646
647
648
649
649
650
651
652
653
654
655
656
657
658
659
659
660
661
662
663
664
665
666
667
668
669
669
670
671
672
673
674
675
676
677
678
679
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
697
698
699
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
897
898
899
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
988
989
989
990
991
992
993
994
995
995
996
997
997
998
999
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1088
1089
1089
1090
1091
1092
1093
1094
1095
1095
1096
1097
1097
1098
1099
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1188
1189
1189
1190
1191
1192
1193
1194
1195
1195
1196
1197
1197
1198
1199
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1288
1289
1289
1290
1291
1292
1293
1294
1295
1295
1296
1297
1297
1298
1299
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1388
1389
1389
1390
1391
1392
1393
1394
1395
1395
1396
1397
1397
1398
1399
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1488
1489
1489
1490
1491
1492
1493
1494
1495
1495
1496
1497
1497
1498
1499
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1588
1589
1589
1590
1591
1592
1593
1594
1595
1595
1596
1597
1597
1598
1599
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1688
1689
1689
1690
1691
1692
1693
1694
1695
1695
1696
1697
1697
1698
1699
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1788
1789
1789
1790
1791
1792
1793
1794
1795
1795
1796
1797
1797
1798
1799
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1849